



Balanced parentheses in NL texts : a useful cue in the syntax/semantics interface

Gabriel G. Bès, Veronica Dahl

► To cite this version:

Gabriel G. Bès, Veronica Dahl. Balanced parentheses in NL texts : a useful cue in the syntax/semantics interface. Workshop on Prospects and Advances in the Syntax/Semantics Interface, 2003, Nancy, France. hal-01096737

HAL Id: hal-01096737

<https://hal.science/hal-01096737>

Submitted on 18 Dec 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Balanced parentheses in NL texts: a useful cue in the syntax/semantics interface*

Gabriel G. Bès[†]
Université Blaise-Pascal
GRIL

Veronica Dahl[‡]
Simon Fraser University
Computing Sciences Department

Abstract

Balanced parentheses on text sentences can be obtained from information on particular morphemes – the introducers – and on inflected verbal forms. From balanced parentheses, a partial graph of the sentence in the semantics interface can be deduced, along with other information. The hypothesis and its expression with CHR constraints are presented.

1 The basic hypothesis

Many formal languages use parentheses, left ones (*lp*) and right ones (*rp*). They are balanced: at the end of a well formed expression $N(lp) = N(rp)$ (where N : number), and, at any point of it, $N(lp) \leq N(rp)$.

Parentheses are well identified objects in formal languages. Classified under the label of "auxiliary symbols" they do not have intrinsic semantic value, but they are crucially important for the specification of operators domains. In Montague Grammar (Montague, 1974), their expressive power is even greater: indexed parentheses encode the syntactic operation from which they follow.

Parentheses are widely used in formal or quasi formal syntactic representations. But, as pointed out by Hintikka (Hintikka, 1994), they are not "natural" objects. They belong to the syntactic machinery of the metalanguage used to describe NL expressions, and as such, their

use can freely change from one machinery to another, even if they remain balanced.

But not all formal languages need parentheses as auxiliary symbols. The polish notation of first order logic does not require them. Our central hypothesis is a kind of an answer to Hintikka's challenge objection: balanced parentheses can indeed be deduced – and not stipulated – from an adequate analysis of NL expressions. Furthermore, they lead to a partial graph, which can be used as an important cue in the syntax/semantic interface.

Balanced parentheses can be obtained from an adequate analysis of a subset of grammatical morphemes such as French *si*, *que*, ..., the introducers, and inflected verbal forms, inflected chunks (e.g. *a lu*, *lui a donné*, or inflected verbs (e.g. *aimait*, *parle*). Metaphorically, they allow to jump to the roof of a sentence from poor information on local marks in its foundations.

2 From local information to the partial graph

The balanced parentheses hypothesis can be illustrated by the following (i), analyzed by the subsequent (ii) to (v).

- i Si les parents s'étaient mis d'accord hier et avaient bien connu la réglementation, les bureaucrates à qui ils se sont adressés aujourd'hui ne leur auraient pas répondu que c'était impossible, ils auraient dû présenter leur dossier autrement.

With respect to (i), it is possible to say that there are inflected nuclear verbal phrases (*vnfl*), as *se sont adressés*, that in one case, two *vnfls* coordinate (*s'étaient mis d'accord* and *avaient bien connu*), that this verbal coordination is the verbal form of the conditional sentence, that the verb form of the root sentence is *ne leur*

Authors in alphabetical order. Thanks are given to Caroline Hagège for extended and enlightening discussions in the preliminaries of this work, and to François Trouilleux for his comments.

Gabriel.Bes@univ-bpclermont.fr; 34 Ave. Carnot, F 63037 Clermont-Fd cedex.

†veronica@sfu.ca; 8888 University Dr. Burnaby B.C. V5A 1S6 Canada.

auraient pas répondu, that the whole sentence with *auraient dû* as verbal form is coordinated to the root sentence.

Futhermore, it is possible to say that there are morphological expressions, simple (as *si*) or complex (as *à qui*), which flag a coming *vnfl*; these are the *introducers*. For instance *se sont adressés* is introduced by *à qui*, *auraient dû* by the nominative form *ils*. The first *vnfl* of the coordinated verbal form of the conditional sentence is introduced by *si* and the *vnfl* of the root sentence is introduced by a hidden *il* (*initial limit*), assumed as first element of any expression, as *fp* (*final point*) is the final one.

If we introduce an *lp* at the left of *il* and an *rp* at the right of *fp*, associate an *lp* to each introducer and an *rp* to each introduced *vnfl* not coordinated with another *vnfl*, and if we associate an *lp* to the first *vnfl* of coordinated *vnfls* and two *rp* to the right of the last coordinated *vnfl*, balanced parentheses on the whole sentence are obtained.

Thus, from (i), the following (ii) is obtained. In (ii), '-' joins single expressions in (i), obtaining chunk expressions which are computed, with respect to position and tags, as the simple ones. A position is assigned in (ii) to each expression jointly with a tag from a very restricted vocabulary $V = \{int, v, v1, v2, il, fp, ot\}$.

Besides *il* and *fp*, presented earlier, *int* in V is associated to introducers, *v* is associated to *vnfls* not immediately preceded either by ',' or by a coordination form, *v1* is associated to *vnfls* immediately preceded by ',', *v2* is associated to *vnfls* immediately preceded by a coordination form (as *et*, *ou...*) not preceded by a ',', and *ot* (*other*) is associated to any expression which is not associated to one of the previous tags. *INT* will spell both *int* and *il*.

ii ((il<0,il> (Si<1,int> les<2,ot> parents<3,ot> (s'-étaient-mis-d'accord<4,v> hier<5,ot> et-avaient-bien-connu<6,v2>)) la<7,ot> réglementation<8,ot> ,<9,ot> les<10,ot> bureaucrates<11,ot> (à-qui<12,int> ils<13,ot> se-sont-adressés<14,v>) aujourd'hui<15,ot> ne-leur-auraient-pas-répondu<16,v>) (que<17,int> la<18,ot> chose<19,ot> était<20,v>) impossible<21,ot> ,<22,ot> (ils<23,int> auraient-dû<24,v>) présenter<25,ot> leur<26,ot> dossier<27,ot> autrement<28,ot> pf)

If we eliminate NL expressions, leaving only parentheses, tags from V and positions, we ob-

tain the more perspicuous (iii) or (iv). In (iv), '...', spelling *intervals*, substitutes for *ots*.

iii ((int₀ (int₁ ot₂ ot₃ (v₄ ot₅ v₂₆)) ot₇ ot₈ ot₉ ot₁₀ ot₁₁ (int₁₂ ot₁₃ v₁₄) ot₁₅ v₁₆) (int₁₇ ot₁₈ ot₁₉ v₂₀) ot₂₁ ot₂₂ (int₂₃ v₂₄) ot₂₅ ot₂₆ ot₂₇ ot₂₈ fp)

iv ((int₀ (int₁... (v₄... v₂₆))... (int₁₂... v₁₄) ... v₁₆)(int₁₇... v₂₀) ... (int₂₃ v₂₄)... fp)

In (iv), besides intervals, there are *INT*'s and *vnfls* (i.e. *v*, *v2*) each in some position *0* to *n*. These relations can be expressed by pairs $\langle i, j \rangle$, where $i \neq j$ and $i, j \geq 0$. These pairs will be assigned to different sets.

By general convention, *q* is the position of the *vnfl* introduced by some *INT*, i.e. the non coordinated *vnfl* associated to the ')' which closes the associated '(', or the first *vnfl* in a coordinated chain of *vnfls*, coordinated chain which closes the *INT*. If *INT* = *il*, we express the relation by $\langle q, 0 \rangle$, if *INT* \neq *il* and in position *p*, by $\langle p, q \rangle$.

Closing pairs (both $\langle q, 0 \rangle$ and $\langle p, q \rangle$) are in the set $Cl[osing]$. From $\langle q, 0 \rangle \in Cl$ we can deduce that $q \in R[oot]$, where R is either an empty set (see §3.3) or a singleton set with the position of the root *vnfl* as member.

Coordinated *vnfls* in verbal phrases, which are in chains *vnfl*_{w1} ... *vnfl*_{wn}, are denoted by coordination pairs $\langle w1, wi \rangle$, where $wi \neq w1$. Coordination pairs of verbal phrases are in the set $C-sv$. A *vnfl* in position *q* and closing some *int* \neq *il*, can be the verbal form of a sentence coordinated to the root sentence (e.g. *auraient-dû*_{<24,v>} in (ii)). In this case, we write $\langle q, 0 \rangle$ and the pair belongs to the $C-r$ set. With these conventions, from (iv) we obtain (v).

(v) $Cl = \{\langle 16, 0 \rangle, \langle 1, 4 \rangle, \langle 12, 14 \rangle, \langle 17, 20 \rangle, \langle 23, 24 \rangle\}$
 $C-sv = \{\langle 6, 4 \rangle\}$
 $R = \{16\}$
 $C-r = \{\langle 24, 0 \rangle\}$

Cl , $C-sv$ and $C-r$ being sets of pairs, from the union of them it is possible to deduce a partial graph, positions in the input being its vertices.

The parsing system that obtains the elements in (v), given the input expression in (i), is organized in Modules I and II. Module I, succinctly presented here, has as input a chain of Ascii codes of NL texts, associated to one or more

sentences, and obtains representations as in (ii), with one or more segmented and enumerated sentences. An interface obtains (iii) from (ii). The challenge of Module I is the disambiguation of expressions such as *si* or *la juge* which can be or not *ints* or *vnfls*, respectively. It is obtained by exploring local contexts. Module II is expressed in two different ways. There is an algorithm (*Algof-c*) which from (iii) obtains (v). The other way is a plain declarative one, making use of CHR constraints.

3 CHR constraints

CHR (Frühwirth and Abdennadher, 2003) is a very powerful multiset rewriting language. Constraints, viewed as pieces of partial information, are formalized as distinguished, predefined predicates in first-order predicate logic.

A constraint program successively generates constraints as it runs, until a solution is found to the problem or no more constraints can be generated. Rules describe how to generate new constraints from those already generated. For instance, we can view symbols in a grammar as constraints upon word boundaries in an input string. Thus *an interesting morning* could be parsed by CHR rules such as:

- (1) `an(X,Y) ==> det(X,Y).`
- (2) `interesting(X,Y) ==> adj(X,Y).`
- (3) `morning(X,Y) ==> noun(X,Y).`
- (4) `start ==> an(1,2),`
`interesting(2,3), morning(3,4).`
- (5) `det(X,Y), adj(Y,Z), noun(Z,W)`
`==> np(X,W).`

The contiguity and order of the *det*, *adj* and *noun* are ensured by the word boundaries; e.g., the *det* ends where the *adj* starts, at point *Y*. *GC* is the set of all the generated constraints derivable from the input string defined through (1) to (4). It will be generated upon the query: `?- start`. In *GC*, we can then select those that solve the problem we are interested in (in this case `np(1,4)`, which tell us our string analyses into a *np*).

3.1 CHR constraints and grammar rules

CHR rules can directly mirror grammar rules: (5) in §3 mirrors `np → det adj noun`, assuming contiguity between `det adj noun`. CHR rules generate constraints, bottom up and left to right, implementing grammar rules.

Given a structure $\dots vnfl_i \dots vnfl_j \dots$, not all *vnfl* (i.e. *v*, *v1*, *v2*) can instantiate *vnfl_i* or *vnfl_j*. For instance, assuming *h* to the left of *i*:

- (1) If $vnfl_i = v2$ and it coordinates with $vnfl_h = v$ or *v1*, then $vnfl_j \neq v2$.
- (2) If $vnfl_i = v2$ and it does not coordinate with $vnfl_h$, or if it coordinates to a $vnfl_h = v2$, then $vnfl_j$ may be a *v2*.

The grammar which mirrors CHR Rules is thus a grammar of type 1 in the Chomsky hierarchy¹.

The format of the input of CHR rules is

`il ... x1 ... xn ... xm ... fp`

where *x_i* is either an *int* or a *vnfl*, and '...' is, here, either an interval or *e(empty)*. The basic challenge of CHR rules is to specify the constraints in *GC* from which arcs in the partial graph can be obtained.

3.2 Arcs from constraints

The whole set *GC* is not needed for obtaining arcs. *GC* is thus the domain of partial functions specifying pairs of the resulting graph in its range. As an illustration on verbal coordination, consider *a regardé, regarde et regardera ce tableau*, with (vi) as its CHR-rules input.

- (vi) `v(1,2), v1(2,3), v2(3,4), ot(4,5),`
`ot(5,6).`

Several grammar rules specify different types of verbal coordination, two of them underlying the specification of the C-sv set related to (vi):

$vOc1 \rightarrow v \ v1$, which coordinates *v1* to *v* obtaining *vOc1*

$vOc \rightarrow v \ v1 \ v2$, which coordinates *v1* and *v2* to *v* obtaining *vOc*

The CHR rules obtain the *GC* (vii) from (vi).

- (vii) `v(1,2), v0(1,2), v1(2,3),`
`v1R(2,3), v0c1(1,3), v1NT(2,3),`
`cv(1,2), v2(3,4), v2R(3,4),`
`v1c(2,4), v0c(1,4), ot(4,5),`
`ot(5,6), otR(4,6), cv(2,5),`
`cv(1,5), ? yes`

All the informations in (vii) are not needed for obtaining elements in *C-sv*. For instance, intervals, expressed by *otR*, are not significant for

¹The grammar and the CHR-rules program can be provided on demand.

the extraction of the partial graph. Among the partial functions with GC as domain, we have $F1$ and $F2$. Given (vii), $\langle 2,1 \rangle, \langle 3,1 \rangle \in C\text{-sv}$ are obtained by $F1$ and $F2$, respectively.

$F1 : v0c1(X,Z), v(X,Y), v1R(Y,Z) \in GC$
 $\rightarrow \langle (Z-1), X \rangle \in C\text{-sv}$

$F2 : v0c(X,W), v0c1(X,Z), v2(Z,W) \in GC$
 $\rightarrow \langle Z, X \rangle \in C\text{-sv}$

Another example illustrates the obtention of the Cl set. Consider the embedded sentences (*dit*) $que_1 la_2 fille_3 que_4 Jacques_5 a\text{-regardée}_6 est\text{-partie}_7$ with (viii) as its input to CHR rules.

(viii) $int(1,2), \dots, int(4,5), \dots, v(6,7),$
 $\dots, v(7,8)$

(1) and (2) in (ix) compact several grammar rules. In (1), X is *cf*, (*constituant fermé*), or *e*, and *iNT* rewrites as *int ot**, *ot** expressing intervals or *e*. In (2), Y explicit contextual restrictions (see §3.1), while Z rewrites *cv* (*complexe verbal*), obtaining terminal strings (a *vnfl* or a *vnfl* coordination) with an initial *vnfl*.

(ix) (1) $cf \rightarrow iNT X cv$
(2) $cv/Y \rightarrow Z ot^*$

There are CHR rules which mirror (ix). In (x), (1) is a subset of the GC obtained by them, (2) is the partial function $F3$. Given (viii), (1) and (2), (3) is obtained.

(x) (1) $\{cf(x,y), cv(y,z)\}$
(2) $F3 : cf(X,Y), cv(Y,Z) \in GC$
 $\rightarrow \langle X, Y \rangle \in Cl$
(3) $\langle 4, 6 \rangle, \langle 1, 7 \rangle \in Cl$

3.3 Deduced information on intervals

From arcs obtained by partial functions from GC , besides the possibility of expressing the semantic representation of verbal coordination, other interesting informations can be deduced. For instance, if $\langle 0, p \rangle \notin Cl$, then $R = \{\}$, and it is likely that the expression will be a nominal phrase, even with one or more embedded relatives. Furthermore, deduced parentheses associated to particular input symbols specify intervals with inherent restrictions, as in (ix).

(xi) (1) $(int_i \dots (int_j$
(2) $(int_i \dots VFORM, \text{ where } VFORM \text{ is } vnfl_j) \text{ or } ((vnfl_j$
(3) $vnfl_i \dots vnfl_j$

A nominal phrase in (1) can be the subject of a *vnfl* in some position to the right of position j , while that is not the case with (2) or (3).

4 Ongoing work and discussion

Ongoing work relates to the extension of verbal coordinations, to the improvement of the expressive power of the grammar (today with less expressive power than *Algof-c*, see §3) and its CHR implementation and, last but not the least, to the evaluation in effective texts of the underlying linguistic hypothesis, knowing beforehand that neither all *ints* or all *vnfls* can be obtained by Module I, nor all verbal coordination be handled by Module II.

Even if we know this, we claim that, in general, from poor and local information obtained by Module I, thanks to the deductible character of balanced parentheses in NL texts, it is possible to obtain a partial graph of the whole sentence, with good and effective approximations in the NL-software engineering domain. From this, besides verbal coordination, it is possible to deduce, in turn, restrictions on intervals, which will reduce the parsing research space².

References

- Luísa Coheur, Nuno Mamede, and Gabriel G. Bès. 2003. Asdecopas: a syntactic-semantic interface. In *EPIA'03 Workshop on Natural Language and Text Retrieval*, Evora (Portugal).
- T. Frühwirth and S. Abdennadher. 2003. *Essentials of Constraint Programming*. Springer Verlag.
- Jaako Hintikka and Gabriel Sandu. 1997. Game theoretical semantics. In van Johan Benthem and Alice ter Meulen, editors, *Handbook of Logic and Language*, pages 361–410. Elsevier.
- Jaako Hintikka. 1994. *Fondements d'une théorie du langage*. PUF.
- Richard Montague. 1974. Universal grammar. In Richard Thomason, editor, *Formal Philosophy, selected papers of Richard Montague*, pages 222–246. Yale University Press.

²The conjecture is that the analysis of intervals will obtain new arcs and vertices, and from them an oriented graph on the sentence, from which, in turn, semantic functions will obtain the semantic representation, cf. (Coheur et al., 2003), but, following (Hintikka, 1994) and (Hintikka and Sandu, 1997), we do not claim that from balanced parentheses scopes of quantifiers can be deduced.